

Documentation of ChatGPT Support for Advanced BA Project

Disclaimer: This report was also generated using Generative AI to collect prompts used in the project.

Group 25

April 25, 2025

Contents

1	Git/GitHub Issues & CI/CD Setup	3
2	Model Design, Architecture, and Evaluation	6
3	Loss Functions and Optimizer Setup	10
4	Model Testing, Unit Tests, and Coverage	13
5	Dataset Loading, Indexing, and Augmentation	17
6	Image Normalization and Prediction on External API Images	20
7	Training, Hyperparameter Tuning, and Evaluation Metrics	22
8	General Debugging and Additional Helpful Tips	25
9	Visualization Techniques	28
10	Dataset Loading	33
11	Common Debugging Issues	34
12	Git/GitHub and CI/CD Setup	35

13 Citations and References	36
14 Baseline Segmentation Model	36
15 Evaluation of the Baseline	37
16 Visualization Function	38
17 Finding High-Forest Samples	39
18 Basic Augmentation Strategies	40
19 API CHAT history	41

1 Git/GitHub Issues & CI/CD Setup

Query

“Why does this not work in .gitignore when I specify `models/` to ignore my model weights?”*

Answer

Git does not automatically stop tracking files that were already committed, even if they are later added to `.gitignore`. To stop tracking a file that has already been committed, use:

```
1 git rm --cached models/your_model_file.pth
2 git commit -m "Remove large model file from repo"
```

This removes the file from version control without deleting it locally. The `.gitignore` then prevents future commits of this file.

Query

“I get the error ‘No module named `advanced_ba_project`’ when running my tests in GitHub Actions. What should I do?”

Answer

The project must be installed as a package in the GitHub Actions environment. Use:

```
1 pip install -e .
```

This allows Python to recognize the `advanced_ba_project` module using the `pyproject.toml` configuration.

Query

“How should my GitHub Actions workflow look for running unit tests with coverage?”

Answer

Here is the correct GitHub Actions configuration:

```

1 name: Unit Tests
2
3 on:
4   push:
5     branches: [main]
6   pull_request:
7     branches: [main]
8
9 jobs:
10  test:
11    runs-on: ubuntu-latest
12    strategy:
13      matrix:
14        python-version: ["3.12", "3.11"]
15
16    steps:
17      - uses: actions/checkout@v4
18      - uses: actions/setup-python@v5
19        with:
20          python-version: ${ matrix.python-version }
21          cache: 'pip'
22
23      - name: Install dependencies
24        run: |
25          python -m pip install -U pip setuptools wheel
26          pip install -r requirements.txt
27          pip install -e .
28
29      - name: Run tests with coverage
30        run: |
31          coverage run --source=advanced_ba_project -m pytest
32          tests/
33          coverage report -m

```

Query

"I got an error with `export HYDRA_FULL_ERROR=1` on Windows runners in my GitHub workflow. How do I fix that?"

Answer

The `export` command is Unix-specific and does not work on Windows. To avoid this problem and support all operating systems, set environment variables like this:

```

1 - name: Run tests & coverage
2   env:

```

```

3     HYDRA_FULL_ERROR: 1
4     run: |
5         coverage run --source=advanced_ba_project -m pytest tests
6         /
          coverage report -m

```

This method works across Linux, macOS, and Windows runners.

Query

“How do I avoid counting the test files when measuring coverage?”

Answer

Add the following configuration to your `pyproject.toml`:

```

1 [tool.coverage.run]
2 omit = ["tests/*"]

```

This ensures that coverage measurement only includes your project code, not the test files themselves.

Query

“I can’t push my model file to GitHub because it exceeds 100MB. What should I do?”

Answer

GitHub has a 100MB limit for individual files. To avoid this:

1. Remove the large file from Git tracking:

```

1     git rm --cached models/your_model_file.pth
2     git commit -m "Remove large model file"
3

```

2. Add the file pattern to your `.gitignore`:

```

1     models/*.pth
2

```

3. Optionally, use Git LFS (Large File Storage): <https://git-lfs.github.com/>

Query

“How can I cache my Python dependencies in GitHub Actions?”

Answer

In your workflow, use the cache option provided by `actions/setup-python`:

```
1 - uses: actions/setup-python@v5
2   with:
3     python-version: ${{ matrix.python-version }}
4     cache: 'pip'
```

This speeds up your CI runs by caching installed dependencies between builds.

2 Model Design, Architecture, and Evaluation

Query

“Can you help me understand the design of my U-Net model and why it looks like this?”

Answer

Your U-Net model follows an encoder-decoder architecture with skip connections between corresponding encoder and decoder layers. This design allows the model to:

- Downsample the input image to extract features (encoder).
- Upsample to reconstruct the segmentation mask (decoder).
- Use skip connections to combine low-level spatial information with high-level features, improving accuracy.

Structure Overview:

- Contracting blocks with convolution, batch normalization, ReLU activation, and dropout.
- Bottleneck layer at the deepest part of the network.
- Expanding blocks with upsampling and concatenation of encoder features.
- Final 1x1 convolution to reduce channels to the output class.

Query

“How does the forward pass in my U-Net work?”

Answer

The forward pass performs:

1. Encoding the input through several downsampling layers.
2. Processing through the bottleneck layer.
3. Decoding through upsampling, concatenating with encoder features (skip connections).
4. Final 1x1 convolution to output the prediction mask.

Query

“How should I describe the skip connections in U-Net?”

Answer

Skip connections directly link each encoder layer to its corresponding decoder layer. This helps:

- Preserve spatial information lost during downsampling.
- Improve localization performance.
- Allow gradients to flow more easily during backpropagation.

Query

“How do I explain the choices for dropout and batch normalization in my model?”

Answer

Batch Normalization helps stabilize training and allows higher learning rates by normalizing layer inputs. **Dropout** adds regularization by randomly dropping units during training, reducing the chance of overfitting.

In your U-Net:

- Batch normalization is applied after convolutions.
- Dropout is used in the contracting blocks and bottleneck, preventing overfitting.

Query

“What input and output sizes should my U-Net handle?”

Answer

The model is designed to handle inputs of shape:

```
1 (batch_size, 3, img_dim, img_dim)
```

and produces outputs of:

```
1 (batch_size, 1, img_dim, img_dim)
```

matching the size of the segmentation masks.

Query

“Can you explain why we use a 1x1 convolution in the final layer of the U-Net?”

Answer

The 1x1 convolution reduces the number of channels from the last decoder layer to the desired output channels (e.g., 1 for binary segmentation). It:

- Acts as a linear classifier at each pixel location.
- Maintains spatial resolution.
- Is computationally efficient.

Query

“What activation function should I use in the output layer?”

Answer

Since your task is binary segmentation:

- Use `sigmoid` activation on the output.
- Follow with a threshold (e.g., 0.5) to obtain binary masks.

Query

“What loss function should I use for segmentation tasks like this?”

Answer

The recommended approach is a combination of:

- Binary Cross-Entropy (BCE) loss: Stable pixel-wise classification.
- Dice loss: Focuses on overlap between predicted masks and ground truth.

The combination balances pixel-level accuracy and global overlap quality.

Query

“How do I evaluate the performance of my segmentation model?”

Answer

Common evaluation metrics for segmentation include:

- Accuracy
- Precision
- Recall
- F1 Score
- Intersection over Union (IoU)

You can also visualize:

- Training and validation loss curves.
- Example predictions vs. ground truth masks.

Query

“Can you give me an example of how to visualize the model predictions?”

Answer

Here's a minimal example for plotting predictions:

```
1 import matplotlib.pyplot as plt
2 import torch
3
4 def visualize_prediction(model, image_tensor):
5     model.eval()
6     with torch.no_grad():
7         pred = model(image_tensor.unsqueeze(0))
8         pred = torch.sigmoid(pred).squeeze().cpu().numpy()
9     plt.imshow(pred > 0.5, cmap='gray')
10    plt.title('Predicted Mask')
11    plt.axis('off')
12    plt.show()
```

This helps qualitatively assess your model's performance.

3 Loss Functions and Optimizer Setup

Query

“What loss function should I use for my segmentation model?”

Answer

The project uses a **combined Dice + Binary Cross-Entropy (BCE) loss**, often referred to as `DiceBCELoss`. This combination was chosen because:

- BCE loss focuses on pixel-wise classification (each pixel as forest/non-forest).
- Dice loss emphasizes overlap between predicted and true masks, which helps especially when there is class imbalance (small forest regions vs. large background).

Query

“How is the Dice + BCE loss implemented?”

Answer

The total loss is calculated as:

$$\text{Loss} = \text{BCE Loss} + \text{Dice Loss}$$

The Dice loss formula is:

$$\text{Dice Loss} = 1 - \frac{2 \cdot (\text{intersection}) + \epsilon}{(\text{sum of predictions}) + (\text{sum of targets}) + \epsilon}$$

where ϵ is a small constant (e.g., $1e-6$) to avoid division by zero.

Query

“Why use both BCE and Dice loss together?”

Answer

BCE loss alone can struggle with imbalanced classes where the background dominates. Dice loss compensates for this by focusing on the overlap area between predictions and targets. Combining the two leverages:

- Stable convergence from BCE.
- Overlap sensitivity from Dice.

Query

“How do I set up the optimizer in PyTorch for my model?”

Answer

The optimizer used in the project is **Adam**. This optimizer is well-suited for segmentation tasks due to its adaptive learning rates. Additionally, **weight_decay** (L2 regularization) is applied to help avoid overfitting.

Example setup:

```
1 optimizer = torch.optim.Adam(  
2     model.parameters(),  
3     lr=cfg.hyperparameters.learning_rate,  
4     weight_decay=cfg.hyperparameters.weight_decay  
5 )
```

Query

“Should I use a learning rate scheduler for this model?”

Answer

Yes, the project uses `ReduceLROnPlateau`. This scheduler:

- Monitors the validation loss.
- Reduces the learning rate by a factor (e.g., 0.5) if no improvement is seen after a certain number of epochs (patience).

Example setup:

```
1 scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(  
2     optimizer,  
3     mode='min',  
4     factor=0.5,  
5     patience=3,  
6     verbose=True  
7 )
```

Query

“Can you show me the full setup for the loss function, optimizer, and scheduler?”

Answer

Here’s the complete setup used in your project:

```
1 criterion = DiceBCELoss()  
2  
3 optimizer = torch.optim.Adam(  
4     model.parameters(),  
5     lr=cfg.hyperparameters.learning_rate,  
6     weight_decay=cfg.hyperparameters.weight_decay  
7 )  
8  
9 scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(  
10     optimizer,  
11     mode='min',  
12     factor=0.5,  
13     patience=3,  
14     verbose=True  
15 )
```

Query

“What hyperparameters should I tune for the optimizer and scheduler?”

Answer

Key hyperparameters tuned during your sweeps:

- `learning_rate`: Tested values like `{0.001, 0.0001}`.
- `weight_decay`: `{0.001, 0.0001}`.
- `scheduler factor`: Commonly 0.5.
- `scheduler patience`: Usually set to 3 epochs.

4 Model Testing, Unit Tests, and Coverage

Query

“How can I write a simple test for the UNet model?”

Answer

The basic test for the U-Net model should verify:

- Whether the model can be instantiated without errors.
- Whether the forward pass runs successfully.
- Whether the output shape matches the expected shape.

Example Unit Test Code

```
1 import torch
2 from advanced_ba_project.model import UNet
3
4 def test_unet_forward_pass():
5     model = UNet(in_channels=3, out_channels=1)
6     x = torch.randn((1, 3, 128, 128))
7     out = model(x)
8     assert out.shape == (1, 1, 128, 128), "Output shape is
    incorrect"
```

Query

“How can I improve this test to check the model structure as well?”

Answer

For a more robust check, you can test both the structure and the forward pass:

```
1 def test_unet_structure_and_forward():
2     model = UNet(in_channels=3, out_channels=1)
3     assert model.in_channels == 3
4     assert model.out_channels == 1
5     assert hasattr(model, 'encoder'), "Encoder not found in
model"
6     assert hasattr(model, 'decoder'), "Decoder not found in
model"
7
8     x = torch.randn((1, 3, 256, 256))
9     out = model(x)
10    assert out.shape == (1, 1, 256, 256), "Incorrect output
shape"
```

Query

“Why do I get an import error: ‘No module named advanced_ba_project’ during testing?”

Answer

This happens if your project is not properly installed as a Python package. The solution is to include:

```
1 pip install -e .
```

This makes `advanced_ba_project` discoverable by Python and your test runner.

Query

“How do I run my tests with coverage tracking?”

Answer

You can use `pytest` together with `coverage.py`:

```
1 coverage run --source=advanced_ba_project -m pytest tests/
2 coverage report -m
```

Query

“How should my GitHub Actions workflow be configured to run these tests?”

Answer

Your CI workflow should:

- Install all dependencies via `requirements.txt`.
- Install your package using `pip install -e ..`
- Run coverage with `pytest`.

Example GitHub Actions Workflow

```
1 name: Unit Tests
2
3 on:
4   push:
5     branches: [main]
6   pull_request:
7     branches: [main]
8
9 jobs:
10  test:
11    runs-on: ubuntu-latest
12    strategy:
13      matrix:
14        python-version: ["3.12", "3.11"]
15
16    steps:
17      - uses: actions/checkout@v4
18      - uses: actions/setup-python@v5
19        with:
20          python-version: ${ matrix.python-version }
21          cache: 'pip'
22
23      - name: Install dependencies
24        run: |
25          python -m pip install -U pip setuptools wheel
26          pip install -r requirements.txt
27          pip install -e .
```

```

28
29     - name: Run tests with coverage
30       run: |
31         coverage run --source=advanced_ba_project -m pytest
        tests/
32         coverage report -m

```

Query

*“Why did I get errors when using **export** in my GitHub Actions on Windows?”*

Answer

The **export** command works on Linux and macOS but not on Windows. For GitHub Actions to work across all operating systems, use:

```

1 env:
2   HYDRA_FULL_ERROR: 1

```

This is cross-platform and safe on Ubuntu, macOS, and Windows runners.

Query

“How do I prevent coverage from including my test files?”

Answer

In your `pyproject.toml`, add:

```

1 [tool.coverage.run]
2 omit = ["tests/*"]

```

Query

“Is it necessary to test performance metrics in my unit tests?”

Answer

No, unit tests should focus on:

- Correct model initialization.
- Forward pass correctness.

- Output shape validation.

Performance metrics (like accuracy or F1 score) should be evaluated separately in your training and evaluation scripts, not in your unit tests.

5 Dataset Loading, Indexing, and Augmentation

Query

“Why don’t I get the same images when I use the same indices between `visualize_raw_and_predicted` and `visualize_raw_and_baseline`?”

Answer

The issue was that the file ordering in your dataset loader was not guaranteed to be consistent across different runs. When using `Path.glob("*.jpg")` or similar functions, the file ordering depends on the filesystem and is not sorted by default.

Solution

Always explicitly sort your list of image files to make indexing reliable:

```
1 image_paths = sorted(Path(image_dir).glob("*.jpg"))
```

This ensures that index 0 always corresponds to the same image across both functions.

Example Fix for RoboflowTreeDataset

Replace:

```
1 self.samples = self._build_samples()
```

with:

```
1 self.samples = sorted(self._build_samples())
```

Query

“Is there any way I can sort the dataset externally without changing the `RoboflowTreeDataset` class?”

Answer

If you do not want to change the dataset class itself, you can still enforce sorting externally by sorting the `indices` list when calling your visualization functions:

```
1 indices = sorted(indices)
```

However, it is safer and more reliable to sort the file paths at the source (within the dataset loader).

Query

“What kind of augmentations were used during training?”

Answer

The project applied basic data augmentations through the `Augmentation` class. These included:

- Random horizontal flips (probability configurable).
- Random vertical flips (probability configurable).
- Random rotations by 90°, 180°, or 270° (probability configurable).

Augmentation Class Example

```
1 class Augmentation(Dataset):
2     def __init__(self, dataset, flip_prob=0.5, rotate_prob
3         =0.3):
4         self.dataset = dataset
5         self.flip_prob = flip_prob
6         self.rotate_prob = rotate_prob
7
8     def __len__(self):
9         return len(self.dataset)
10
11     def __getitem__(self, idx):
12         image, mask = self.dataset[idx]
13
14         if random.random() < self.flip_prob:
15             image = torch.flip(image, [2])
16             mask = torch.flip(mask, [2])
17
18         if random.random() < self.flip_prob:
19             image = torch.flip(image, [1])
```

```

19         mask = torch.flip(mask, [1])
20
21     if random.random() < self.rotate_prob:
22         k = random.choice([1, 2, 3])
23         image = torch.rot90(image, k, [1, 2])
24         mask = torch.rot90(mask, k, [1, 2])
25
26     return image, mask

```

Query

“How were the training, validation, and test sets split?”

Answer

The dataset was split into:

- 80% training.
- 10% validation.
- 10% test.

This was done using `torch.utils.data.random_split` with a manual seed for reproducibility.

Query

“How do I ensure that the validation and test sets are consistent between runs?”

Answer

Set the random seed explicitly before splitting:

```

1 torch.manual_seed(seed)
2 train_size = int(0.8 * len(dataset))
3 val_size = int(0.1 * len(dataset))
4 test_size = len(dataset) - train_size - val_size
5 train_set, val_set, test_set = random_split(dataset, [
    train_size, val_size, test_size])

```

This ensures the splits remain consistent as long as the dataset order remains consistent.

Key Takeaways

- Always sort file paths when loading data.
- Apply a fixed seed when splitting datasets to keep validation and test sets stable.
- Use augmentation to improve model generalization.

6 Image Normalization and Prediction on External API Images

Query

“Why does my model work well on the training images but not on the API-provided .tif images?”

Answer

The model was trained on images normalized with:

```
1 transforms.Normalize(mean=[0.5, 0.5, 0.5], std=[0.5, 0.5, 0.5])
```

However, the API images had pixel values outside of the expected range and were not scaled between 0 and 1 before normalization. This caused the model to receive data in a very different format from what it saw during training.

Solution

Ensure that all input images are first scaled to the [0, 1] range before applying normalization. The correct approach is:

```
1 img_array = img_array.astype(np.float32)
2 img_array = (img_array - img_array.min()) / (img_array.max()
    - img_array.min() + 1e-8)
```

Query

“Can I simply scale the API images between 0 and 1 and then apply the same normalization?”

Answer

Yes, this is the recommended approach. First scale the images to $[0, 1]$, and then apply the exact same normalization used in your training process to ensure consistency.

Corrected Image Loader Function

```
1 def load_tif_image(path, img_dim=256):
2     img = Image.open(path)
3     img_array = np.array(img)
4
5     # Handle multi-channel images or grayscale
6     if img_array.ndim == 3 and img_array.shape[2] > 3:
7         img_array = img_array[:, :, :3]
8     if img_array.ndim == 2:
9         img_array = np.stack([img_array] * 3, axis=-1)
10
11     # Scale to [0, 1]
12     img_array = img_array.astype(np.float32)
13     img_array = (img_array - img_array.min()) / (img_array.
14 max() - img_array.min() + 1e-8)
15
16     img_pil = Image.fromarray((img_array * 255).astype(np.
17 uint8))
18     transform = transforms.Compose([
19         transforms.Resize((img_dim, img_dim)),
20         transforms.ToTensor(),
21         transforms.Normalize(mean=[0.5, 0.5, 0.5], std=[0.5,
22 0.5, 0.5])
23     ])
24     return transform(img_pil)
```

Query

“Why did scaling fix the issue?”

Answer

The neural network expects input values to be normalized around zero with a standard deviation of one. Without scaling, the `.tif` images had unpredictable ranges (sometimes 0–65535, not 0–255), so the normalization step failed to align them properly with the trained data distribution.

Scaling ensures the model sees images in the same format as during training.

Query

“Could histogram matching or equalization work instead of scaling?”

Answer

Histogram matching or equalization could help if the issue was related to contrast or brightness differences between datasets. However, the primary problem here was not contrast—it was the range of pixel values. Simple scaling to $[0, 1]$ is sufficient and more consistent with the training pipeline.

Debugging Insight

Plotting histograms of the pixel distributions from your API images and comparing them to your training set helped identify the scaling mismatch. Without proper scaling, the pixel distributions looked drastically different.

Observed Result

After applying scaling:

- Predictions on API images became significantly more stable.
- Although performance was still not perfect (due to dataset domain differences), the output masks were reasonable.

Additional Recommendation

For future work, consider:

- Including API-style images in your training dataset.
- Using domain adaptation techniques or augmentations to simulate such images during training.

7 Training, Hyperparameter Tuning, and Evaluation Metrics

Query

“How should I describe the hyperparameter tuning and training process?”

Answer

Hyperparameter tuning was performed using **Weights & Biases (wandb)** Sweeps with a grid search strategy. The main metric optimized was the **validation F1 score**.

Hyperparameters Tuned

The sweep explored the following parameter combinations:

- `batch_size`: {16, 32}
- `learning_rate`: {0.001, 0.0001}
- `weight_decay`: {0.001, 0.0001}
- `dropout_rate`: {0.1, 0.2, 0.3}
- `apply_augmentation`: {true, false}

Query

“How should I explain the data augmentation strategy used?”

Answer

Data augmentation was applied during training to improve generalization. The following augmentations were used:

- Random horizontal flips.
- Random vertical flips.
- Random rotations (90°, 180°, or 270°).

These augmentations were applied to both input images and their corresponding masks. The logic is implemented in the **Augmentation** class located in `data.py`.

Query

“How were the datasets split for training, validation, and testing?”

Answer

The project combined data from the **DeepGlobe** and **Roboflow** datasets with the following splits:

- 80% for training.
- 10% for validation.
- 10% for testing.

Query

“What learning rate scheduler was used?”

Answer

The **ReduceLROnPlateau** scheduler was used, monitoring the validation loss. If no improvement was observed for 3 epochs, the learning rate was reduced by a factor of 0.5.

Scheduler Example Code

```
1 scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(  
2     optimizer, mode='min', factor=0.5, patience=3, verbose=  
3     True  
4 )
```

Query

“How were the final results evaluated and reported?”

Answer

Performance evaluation included:

- Validation and test loss.
- Accuracy.
- Precision.
- Recall.

- F1 score.
- IoU (Intersection over Union).

Metrics Logging Example with wandb

```
1 wandb.log({
2     "Validation Loss": val_loss,
3     "Val Accuracy": metrics["accuracy"],
4     "Val Precision": metrics["precision"],
5     "Val Recall": metrics["recall"],
6     "Val F1": metrics["f1"],
7     "Val IoU": metrics["iou"],
8 }, step=epoch + 1)
```

Query

“How were the best hyperparameters determined?”

Answer

The sweep results were analyzed via Weights & Biases to identify the hyperparameter combination that maximized the validation F1 score. This combination was then selected for training the final model.

Observed Results

- Dropout rate and learning rate had significant impact on generalization.
- Applying augmentation generally improved the F1 score.
- Larger batch sizes worked better when GPU memory allowed.

8 General Debugging and Additional Helpful Tips

Query

“How can I get clearer error messages when using Hydra in my GitHub Actions setup?”

Answer

Hydra sometimes suppresses error output by default. To ensure that full stack traces are displayed, set:

```
1 HYDRA_FULL_ERROR=1
```

In GitHub Actions, for cross-platform compatibility (Linux, Windows, MacOS), use:

```
1 env:
2   HYDRA_FULL_ERROR: 1
```

Query

“Why did my CI/CD pipeline break with the error: ‘The term ‘export’ is not recognized as the name of a cmdlet...’?”

Answer

This occurred because the `export` command is specific to Unix-like shells. On Windows, use:

```
1 set HYDRA_FULL_ERROR=1
```

Or, for cross-platform workflows (recommended):

```
1 env:
2   HYDRA_FULL_ERROR: 1
```

Query

“How can I avoid large model files causing GitHub push errors?”

Answer

GitHub has a 100MB file size limit. To avoid this issue:

- Untrack large files already committed:

```
1 git rm --cached path/to/large_file.pth
2 git commit -m "Remove large file from repo"
3
```

- Add these files to `.gitignore` to prevent future commits.
- Consider using **Git LFS (Large File Storage)**:

```
1 git lfs install
2 git lfs track "*.pth"
3 git add .gitattributes
4 git add path/to/large_file.pth
5 git commit -m "Track large model file with Git LFS"
6
```

Query

“What helped you quickly identify the issue with image normalization mismatches?”

Answer

The following debugging steps were effective:

- Plot histograms of pixel values for both training images and API images.
- Compare mean and standard deviation statistics between datasets.
- Visualize example images before and after normalization to confirm the preprocessing pipeline is consistent.

Query

“Any general tips for reproducibility and debugging?”

Answer

Good practices for reproducibility and debugging included:

- Setting random seeds for Python, NumPy, and PyTorch:

```
1 import random
2 import numpy as np
3 import torch
4
5 random.seed(42)
6 np.random.seed(42)
7 torch.manual_seed(42)
8 if torch.cuda.is_available():
9     torch.cuda.manual_seed_all(42)
10
```

- Logging configuration parameters (e.g., via Hydra or wandb).
- Keeping training and evaluation environments consistent.

FROM DOWNWARDS THE PROMPTS ORIGINATE FROM ANOTHER CHAT HISTORY.

9 Visualization Techniques

Query

How can I visualize segmentation masks overlayed on original images?

Answer

Implemented a custom function with color-coded overlays:

```
1 def visualize_batch_from_dataset(dataset, num_images=5,
2   subset_index=None):
3     if subset_index is not None and isinstance(dataset, torch
4       .utils.data.ConcatDataset):
5       dataset = dataset.datasets[subset_index]
6       indices = torch.randperm(len(dataset))[:num_images]
7       images, masks = zip(*[dataset[i] for i in indices])
8       for i in range(num_images):
9         overlay[mask_np==1.0] = [0, 1, 0] # Green
10        overlay[mask_np==0.0] = [0.5, 0, 0] # Dark red
```

Query

How can I choose which dataset to visualize within a ConcatDataset?

Answer

Added a subset_index parameter:

```
1 visualize_batch_from_dataset(val_loader.dataset, num_images
2   =5, subset_index=0) # Forest
3 visualize_batch_from_dataset(val_loader.dataset, num_images
4   =5, subset_index=1) # Roboflow
```

Query

What RGB value represents dark red for mask overlays?

Answer

Dark red is represented by:

```
1 overlay[non_forest_mask] = (139, 0, 0)
```

This color choice supports clarity and accessibility.

Query

How can I visualize segmented forest areas from two different datasets ("forest" and "roboflow") with consistent overlay coloring?

Answer

The following Python function provides a robust method to visualize segmentation masks overlaid on satellite images from either the Forest or Roboflow dataset. It includes options for setting random seeds for reproducibility or specifying exact indices for visualization.

```
1 from advanced_ba_project.data import ForestDataset,
   RoboflowTreeDataset
2 from pathlib import Path
3 import torch
4 import random
5 import numpy as np
6 import matplotlib.pyplot as plt
7 from torchvision import transforms
8 from matplotlib.patches import Patch
9
10 def plot_forest_overlay_grid_from_dataset(
11     dataset_name: str,
12     data_path: Path,
13     metadata_file: str = None,
14     image_dir: Path = None,
15     label_dir: Path = None,
16     num_images: int = 5,
17     img_dim: int = 256,
18     seed: int = 42,
19     indices: list = None
20 ):
21     transform = transforms.Compose([
22         transforms.Resize((img_dim, img_dim)),
23         transforms.ToTensor(),
24         transforms.Normalize(mean=[0.5]*3, std=[0.5]*3)
25     ])
26     target_transform = transforms.Compose([
27         transforms.Resize((img_dim, img_dim)),
```

```

28 transforms.ToTensor()
29 ])
30
31 if dataset_name == "forest":
32     if metadata_file is None:
33         raise ValueError("metadata_file must be provided for
34 forest dataset.")
35     dataset = ForestDataset(
36         data_path=data_path,
37         metadata_file=metadata_file,
38         transform=transform,
39         target_transform=target_transform
40     )
41 elif dataset_name == "roboflow":
42     if image_dir is None or label_dir is None:
43         raise ValueError("image_dir and label_dir must be
44 provided for roboflow dataset.")
45     dataset = RoboflowTreeDataset(
46         image_dir=image_dir,
47         label_dir=label_dir,
48         patch_size=img_dim,
49         transform=transform,
50         target_transform=target_transform
51     )
52 else:
53     raise ValueError("dataset_name must be either 'forest' or
54 'roboflow'.")
55
56 if indices is None:
57     torch.manual_seed(seed)
58     random.seed(seed)
59     indices = random.sample(range(len(dataset)), min(
60 num_images, len(dataset)))
61 else:
62     num_images = len(indices)
63
64 samples = [dataset[i] for i in indices]
65 fig, axs = plt.subplots(2, num_images, figsize=(4 *
66 num_images, 8))
67
68 for i, (image, mask) in enumerate(samples):
69     image = image.permute(1, 2, 0).numpy()
70     image = (image * 0.5) + 0.5
71     image = np.clip(image, 0, 1)
72
73     mask = mask.squeeze().numpy()
74     if mask.ndim == 3:
75         mask = mask[0]

```

```

72     axs[0, i].imshow(image)
73     axs[0, i].axis("off")
74     axs[0, i].set_title(f"{dataset_name.capitalize()}
Original {indices[i]}")
75
76     overlay = image.copy()
77     forest_color = np.array([0, 1, 0])
78     nonforest_color = np.array([1, 0.3, 0])
79     alpha = 0.5
80
81     mask_bool = mask > 0.5
82     overlay[mask_bool] = (1 - alpha) * overlay[mask_bool] +
alpha * forest_color
83     overlay[~mask_bool] = (1 - alpha) * overlay[~mask_bool] +
alpha * nonforest_color
84
85     axs[1, i].imshow(overlay)
86     axs[1, i].axis("off")
87     axs[1, i].set_title(f"{dataset_name.capitalize()} Overlay
{indices[i]}")
88
89 legend_elements = [
90     Patch(facecolor='green', edgecolor='black', label='Forest
', alpha=0.5),
91     Patch(facecolor='orangered', edgecolor='black', label='
Non-Forest', alpha=0.5),
92 ]
93 fig.legend(handles=legend_elements, loc='lower center', ncol
=2, fontsize='large')
94 plt.tight_layout(rect=[0, 0.05, 1, 1])
95 plt.show()

```

Listing 1: Function to plot overlay grids from dataset

Usage Example

```

1 plot_forest_overlay_grid_from_dataset(
2     dataset_name="roboflow",
3     data_path=None,
4     image_dir=Path("data/raw/roboflow/test/images"),
5     label_dir=Path("data/raw/roboflow/test/labelTxt"),
6     seed=100,
7     num_images=5
8 )
9
10 plot_forest_overlay_grid_from_dataset(
11     dataset_name="forest",
12     data_path=Path("data/raw/Forest Segmented"),

```

```

13 metadata_file="meta_data.csv",
14 indices=[0, 3, 7, 12, 25]
15 )

```

Query

How can I ensure that my quick sample visualizations use the same forest/non-forest overlay colors?

Answer

Use this streamlined function to visualize samples consistently across datasets, applying the same green and orange-red overlay scheme as in the main function.

```

1 from matplotlib.patches import Patch
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 def visualize_samples(dataset, num_samples=5, alpha=0.5):
6     fig, axes = plt.subplots(1, num_samples, figsize=(4 *
7         num_samples, 4))
8     for i in range(num_samples):
9         image, mask = dataset[i]
10        image = image.permute(1, 2, 0).numpy()
11        image = (image * 0.5) + 0.5
12        image = np.clip(image, 0, 1)
13
14        mask = mask.squeeze().numpy()
15        if mask.ndim == 3:
16            mask = mask[0]
17
18        overlay = image.copy()
19        forest_color = np.array([0, 1, 0])
20        nonforest_color = np.array([1, 0.3, 0])
21        mask_bool = mask > 0.5
22
23        overlay[mask_bool] = (1 - alpha) * overlay[mask_bool] +
24            alpha * forest_color
25        overlay[~mask_bool] = (1 - alpha) * overlay[~mask_bool] +
26            alpha * nonforest_color
27
28        axes[i].imshow(overlay)
29        axes[i].axis("off")
30        axes[i].set_title(f"Sample {i+1}")
31
32    legend_elements = [

```



```

30     Patch(facecolor='green', edgecolor='black', label='Forest
    ', alpha=alpha),
31     Patch(facecolor='orangered', edgecolor='black', label='
    Non-Forest', alpha=alpha),
32 ]
33 fig.legend(handles=legend_elements, loc='lower center', ncol
    =2, fontsize='large')
34 plt.tight_layout(rect=[0, 0.05, 1, 1])
35 plt.show()

```

Listing 2: Consistent overlay visualization for quick samples

Usage Example

```

1 visualize_samples(roboflow_train, num_samples=5)
2 visualize_samples(forest_dataset, num_samples=3)

```

10 Dataset Loading

Query

How should I organize my dataset for forest segmentation tasks in PyTorch? (image upload of current setup and I also already uploaded example code from already implemented dataloader)

Answer

Two datasets were used:

- **Forest Segmented:** Contains RGB satellite images and float-based binary masks (0.0 for non-forest, 1.0 for forest).
- **Roboflow:** Contains images with polygon annotations (converted into masks).

Folder structure:

```

1 data/
2     raw/
3         Forest Segmented/
4             images/
5             masks/
6             meta_data.csv
7     roboflow/
8         train/
9         valid/

```

Query

How should the dataset classes be structured for loading both datasets?

Answer

Two custom dataset classes were implemented:

ForestDataset:

```
1 class ForestDataset(Dataset):
2     def __init__(self, data_path, metadata_file, transform=
      None, target_transform=None):
3         self.image_dir = data_path / "images"
4         self.mask_dir = data_path / "masks"
5         self.metadata = pd.read_csv(data_path / metadata_file
      )
6
7     def __getitem__(self, index):
8         image = Image.open(self.image_paths[index]).convert("
      RGB")
9         mask = Image.open(self.mask_paths[index]).convert("L"
      )
10        return image, mask
```

11 Common Debugging Issues

Query

Why does my mask overlay not appear as expected? The overlay looks like the original image without any changes.

Answer

This often occurs because the mask is not binary. To ensure proper overlay, you need to explicitly binarize the mask:

```
1 mask_np = mask.squeeze().numpy()
2 binary_mask = (mask_np > 0.5).astype(np.uint8)
```

Listing 3: Binarizing the Mask

Query

Why do I get an `ImportError` when trying to import my module in a notebook?

Answer

Make sure that the `src/` directory is added to `sys.path` at the beginning of your notebook:

```
1 import sys
2 from pathlib import Path
3
4 sys.path.append(str(Path("../src").resolve()))
```

Listing 4: Adding `src` to `sys.path`

This allows proper module imports like:

```
1 from project_name.data import visualize_raw_data
```

12 Git/GitHub and CI/CD Setup

Query

How do I pull everything from the main branch into my branch `explainer-Notebook`?

Answer

Ensure you're on the correct branch, then pull the latest changes from main:

```
1 git checkout explainerNotebook
2 git pull origin main
```

To rebase instead of merging (for a cleaner history):

```
1 git checkout explainerNotebook
2 git fetch origin
3 git rebase origin/main
```

Query

How do I push changes to the `explainerNotebook` branch?

Answer

Follow these steps to push your changes:

```
1 git checkout explainerNotebook
2 git add .
3 git commit -m "Add/update explainer notebook content"
4 git push origin explainerNotebook
```

13 Citations and References

Query

How do I cite Weights & Biases (W&B)?

Answer

APA style:

Weights & Biases. (n.d.). *Weights & Biases Documentation*. Retrieved April 23, 2025, from <https://wandb.ai/site/>
FROM DOWNWARDS THE PROMPTS ORIGINATE FROM ANOTHER CHAT HISTORY.

14 Baseline Segmentation Model

Query

Help me make a baseline for my project that is a segmentation model to detect forest cover in satellite images. The data consists of RGB matrices and a ground truth mask. I want the baseline to check each pixel value and return a 1 in the predicted mask if the pixel value is more green than red and blue. Also include a threshold as input that can be adjusted to say how much larger the green value needs to be in order for it to be predicted as a tree.

Answer

Below is a simple baseline segmentation model implementation based on pixel-wise green dominance:

```
1 def green_tree_detector(image_batch, threshold=0.1):
2     """
3     Simple baseline that detects trees based on green channel
4     values.
5     Args:
6         image_batch: Tensor of shape [batch_size, 3, 256,
7         256]
8         threshold: How much greener a pixel must be compared
9         to other channels
10    Returns:
11        Tensor of shape [batch_size, 1, 256, 256] with binary
12        tree mask
13    """
14    batch_size = image_batch.shape[0]
```

```

11     device = image_batch.device
12
13     r = image_batch[:, 0]
14     g = image_batch[:, 1]
15     b = image_batch[:, 2]
16
17     tree_mask = ((g > (r + threshold)) & (g > (b + threshold)
18     )).float()
19     return tree_mask.unsqueeze(1)

```

15 Evaluation of the Baseline

Query

Evaluate the baseline model using common metrics.

Answer

The following function evaluates the model using accuracy, precision, recall, F1 score, and IoU:

```

1 def evaluate_detector(dataloader, threshold=0.1):
2     device = next(iter(dataloader))[0].device
3
4     metrics = {'accuracy': 0, 'precision': 0, 'recall': 0, '
5     f1': 0, 'iou': 0}
6     total_pixels = 0
7
8     for images, masks in tqdm(dataloader, desc="Evaluating
9     baseline detector"):
10         predictions = green_tree_detector(images, threshold=
11         threshold)
12         pred_binary = (predictions > 0.5).float()
13
14         pred_flat = pred_binary.cpu().numpy().flatten().
15         astype(int)
16         mask_flat = masks.cpu().numpy().flatten().astype(int)
17
18         batch_pixels = pred_flat.shape[0]
19         total_pixels += batch_pixels
20
21         metrics['accuracy'] += accuracy_score(mask_flat,
22         pred_flat) * batch_pixels
23         metrics['precision'] += precision_score(mask_flat,
24         pred_flat, zero_division=0) * batch_pixels
25         metrics['recall'] += recall_score(mask_flat,
26         pred_flat, zero_division=0) * batch_pixels

```

```

20         metrics['f1'] += f1_score(mask_flat, pred_flat,
zero_division=0) * batch_pixels
21         metrics['iou'] += jaccard_score(mask_flat, pred_flat,
zero_division=0) * batch_pixels
22
23     for key in metrics:
24         metrics[key] /= total_pixels
25
26     return metrics

```

16 Visualization Function

Query

Now help me make a visualization function that can visualize the original image, the predicted mask and the original ground truth mask.

Answer

The following function visualizes and saves predictions:

```

1 def visualize_baseline_predictions(val_loader, device,
green_tree_detector, threshold=0.01, num_samples=5):
2     timestamp = datetime.datetime.now().strftime("%Y-%m-%d_%H
-%M-%S")
3     os.makedirs("reports/figures", exist_ok=True)
4     save_path = f"reports/figures/baseline_mask_comparison_{
timestamp}.png"
5
6     images, true_masks = next(iter(val_loader))
7     images, true_masks = images.to(device), true_masks.to(
device)
8
9     with torch.no_grad():
10         predicted_masks = green_tree_detector(images,
threshold=threshold)
11         predicted_masks = (predicted_masks > 0.5).float()
12
13     images = images.cpu().numpy().transpose(0, 2, 3, 1)
14     true_masks = true_masks.cpu().numpy().squeeze(1)
15     predicted_masks = predicted_masks.cpu().numpy().squeeze
(1)
16
17     fig, axes = plt.subplots(num_samples, 3, figsize=(10,
num_samples * 3))
18     for i in range(num_samples):

```

```

19         axes[i, 0].imshow((images[i] * 0.5) + 0.5)
20         axes[i, 0].set_title("Original Image")
21         axes[i, 0].axis("off")
22
23         axes[i, 1].imshow(true_masks[i], cmap="gray")
24         axes[i, 1].set_title("Ground Truth Mask")
25         axes[i, 1].axis("off")
26
27         axes[i, 2].imshow(predicted_masks[i], cmap="gray")
28         axes[i, 2].set_title(f"Baseline Mask (threshold={
threshold:.2f})")
29         axes[i, 2].axis("off")
30
31     plt.tight_layout()
32     plt.savefig(save_path)
33     print(f"[INFO] Baseline mask comparison saved as {
save_path}")
34     plt.show()

```

17 Finding High-Forest Samples

Query

Now help me create a function that takes in our dataloader function and finds the images in the data with the most forest cover in the ground truth mask.

Answer

The following function finds the top-k samples with the most forest pixels:

```

1 def get_top_forest_ids(
2     dataset_name: str,
3     data_path: Path,
4     metadata_file: str = None,
5     image_dir: Path = None,
6     label_dir: Path = None,
7     img_dim: int = 256,
8     top_k: int = 5
9 ):
10     transform = transforms.Compose([
11         transforms.Resize((img_dim, img_dim)),
12         transforms.ToTensor(),
13         transforms.Normalize(mean=[0.5]*3, std=[0.5]*3)
14     ])
15     target_transform = transforms.Compose([
16         transforms.Resize((img_dim, img_dim)),

```

```

17         transforms.ToTensor()
18     ])
19
20     if dataset_name == "forest":
21         dataset = ForestDataset(
22             data_path=data_path,
23             metadata_file=metadata_file,
24             transform=transform,
25             target_transform=target_transform
26         )
27     elif dataset_name == "roboflow":
28         dataset = RoboflowTreeDataset(
29             image_dir=image_dir,
30             label_dir=label_dir,
31             patch_size=img_dim,
32             transform=transform,
33             target_transform=target_transform
34         )
35     else:
36         raise ValueError("dataset_name must be either 'forest'
37                             or 'roboflow'")
38
39     forest_coverage = []
40     for idx in range(len(dataset)):
41         _, mask = dataset[idx]
42         forest_pixels = (mask > 0.5).sum().item()
43         forest_coverage.append((idx, forest_pixels))
44
45     top_ids = sorted(forest_coverage, key=lambda x: x[1],
46                     reverse=True)[:top_k]
47     return [idx for idx, _ in top_ids]

```

18 Basic Augmentation Strategies

Query

I want to augment our data when training, what are general ways to augment an image?

Answer

Here are standard augmentation techniques useful for training:

- Horizontal flip:

```
1 torch.flip(image, [2])
```


- **Vertical flip:**

```
1 torch.flip(image, [1])
```

- **Random rotation** (by 90°, 180°, 270°):

```
1 k = random.choice([1,2,3])  
2 torch.rot90(image, k, [1, 2])
```

19 API CHAT history

<https://chatgpt.com/share/680b6dd4-7094-8004-8763-e74c8cf81a92>